# NCBI's Entrez System

Alex E. Lash, MD
Cancer Imaging Informatics Meeting
Bethesda, MD

# Paris, 1830



Georges
Cuvier
(1769-1832)



Étienne
Geoffroy St. Hilaire
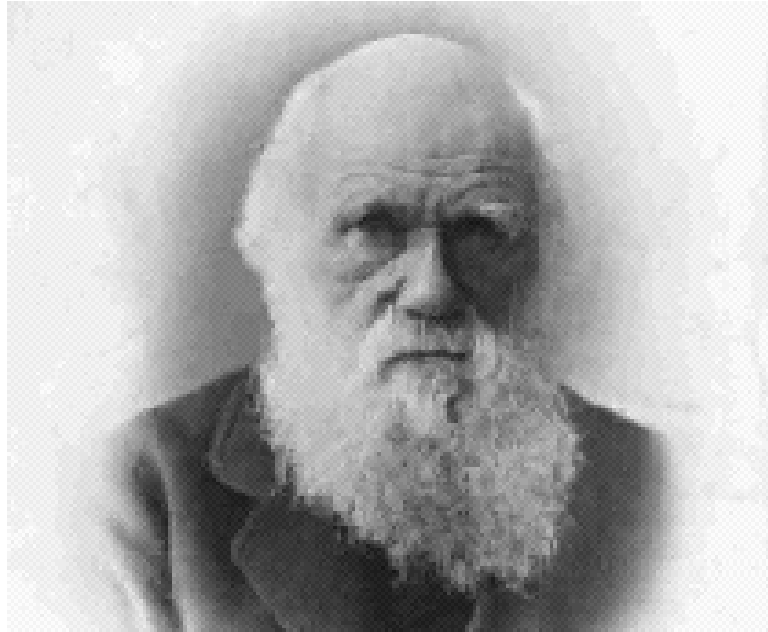(1772-1844)

# 1830: "Form vs. Function" Debate

Cuvier

- "form follows function"
- anatomic similarities among vertebrates were due to similar function
- "If there are resemblances between the organs…, it is only insofar as there are resemblances between their functions."

Geoffroy

- "function follows form"
- vertebrates were modifications of a single archetype
- "There is, philosophically speaking, only a single animal."

# 1859: Darwin on Geoffroy

"Geoffroy St. Hilaire has insisted strongly on the high importance of relative connexion in homologous organs; the parts may change to almost any extent in form and size, and yet they always remain connected together in the same order."

# "Pre-hypothesis" Biological Information Collection

Cuvier & Geoffroy
both got here,
but through different reasoning

⇓

Collect → Characterize → Relate

⇧

where discovery takes place:
patterns are seen
and hypotheses form

A modern
example:

Sequencing:
sequence gene

Annotation:
annotate features
such as coding and
non-coding regions

Cross-comparison:
compare sequence
to every other sequence

# Today vs. 1830

Biotechnological developments have increased size, scope and speed of "pre-hypothesis" biological information collection.

Collection: overwhelming amount and variety of records

GenBank contains >19 million sequence records and >20 billion bases and doubled in size in the last 16 months

Characterization: increased scope and detail of fields in records

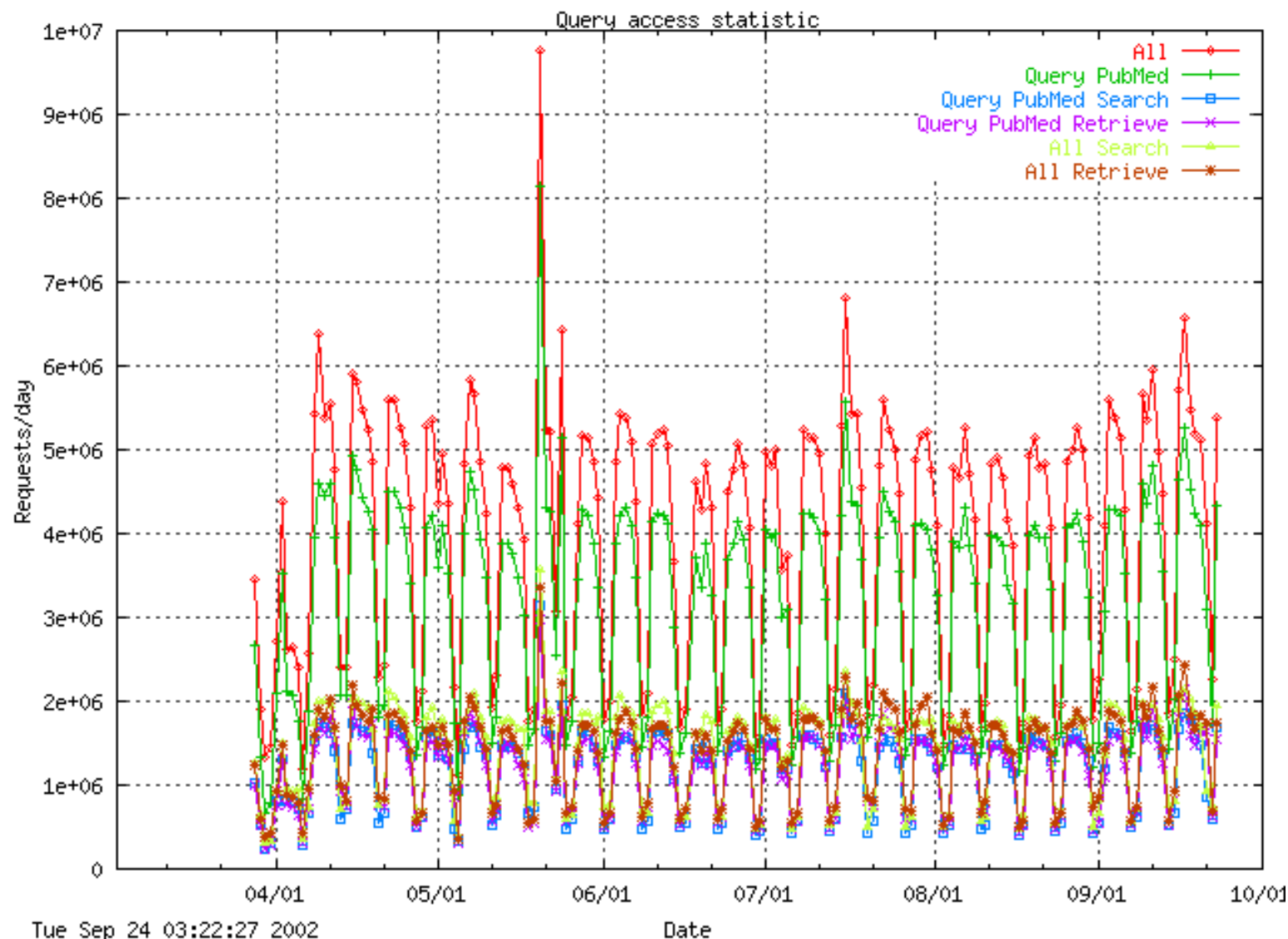Relation: increased possibility of intra- and inter-database record to record links

# National Center for Biotechnology Information

- Created by Public Law 100-607 in 1988 as part of National Library of Medicine at NIH to:

    - Create automated systems for knowledge about molecular biology, biochemistry, and genetics.

    - Perform research into advanced methods of analyzing and interpreting molecular biology data.

    - Enable biotechnology researchers and medical care personnel to use the systems and methods developed.

- Builders and providers of GenBank, Entrez, Blast, PubMed. Online systems host more than 2 million users per month.

- Center for basic research and training in computational biology.

# NCBI Web Hits Per Day



NCBI Web Site

Jan/02

Jan/01

Jan/00

Jan/99

Days since Jan. 1, 1998

# Entrez Hits Per Day



Query access statistic

# What is Entrez?

Entrez is a <u>scalable</u> and <u>flexible</u> database and interface system constructed and maintained at NCBI.

Each Entrez database contains records with pre-specified fields, contains indices on each field, and comes with an interface allowing field-specific, boolean queries.

PubMed is an Entrez database.  OMIM is an Entrez database. GenBank nucleotide sequence records are contained in Entrez Nucleotide.

Links can be specified between records within the same Entrez database (intra-database links), or between records in different Entrez databases (inter-database links).

Links can be obvious (eg, identifier matching) or non-obvious (eg, sequence similarity).  Non-obvious links generally require examination of the full record and some computation.

# Architecture

# Entrez stats

NCBI

15 Entrez databases

>38 million records

>140 million indexed terms

>6.7 billion intra- and inter-database links

# Using Entrez for Discovery - 1

# Using Entrez for Discovery - 2

# Using Entrez for Discovery - 3

# Using Entrez for Discovery - 4

# Using Entrez for Discovery - 5

# Using Entrez for Discovery - 6



NCBI Sequence Viewer - Netscape 6

File  Edit  View  Search  Go  Bookmarks  Tasks  Help

http://www.ncbi.nlm.nih.gov/entrez/viewer.fcgi?val=6323047

Home  Bookmarks  geo castor-po...  geo yar public  geo yar private  SAGEmap  Google  Digital City: W...

Protein

| PubMed | Nucleotide | Protein | Genome | Structure | PopSet | Taxonomy | OMIM | Books |

Search Nucleotide for [ ]  Go  Clear

Limits        Preview/Index        History        Clipboard        Details

Display default  Save  Text  Add to Clipboard  Get Subsequence

```
LOCUS           NP_013119                397 aa            linear   PLN 09-SEP-2002
DEFINITION      Plasma membrane Sodium Response 2; Psr2p [Saccharomyces
                cerevisiae].
DBSOURCE    REFSEQ: accession NC_001144.2
KEYWORDS    .
SOURCE      baker's yeast.
  ORGANISM  Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycetes; Saccharomycetales;
            Saccharomycetaceae; Saccharomyces.
REFERENCE   1  (residues 1 to 397)
  AUTHORS   Goffeau,A., Barrell,B.G., Bussey,H., Davis,R.W., Dujon,B.,
            Feldmann,H., Galibert,F., Hoheisel,J.D., Jacq,C., Johnston,M.,
            Louis,E.J., Mewes,H.W., Murakami,Y., Philippsen,P., Tettelin,H. and
            Oliver,S.G.
  TITLE     Life with 6000 genes
  JOURNAL   Science 274 (5287), 546 (1996)
  MEDLINE   97002444
REFERENCE   2  (residues 1 to 397)
  AUTHORS   Johnston,M., Hillier,L., Riles,L., Albermann,K., Andre,B.,
            Ansorge,W., Benes,V., Bruckner,M., Delius,H., Dubois,E.,
            Dusterhoft,A., Entian,K.D., Floeth,M., Goffeau,A., Hebling,U.,
            Heumann,K., Heuss-Neitzel,D., Hilbert,H., Hilger,F., Kleine,K.,
            Kotter,P., Louis,E.J., Messenguy,F., Mewes,H.W., Hoheisel,J.D. et
            al.
  TITLE     The nucleotide sequence of Saccharomyces cerevisiae chromosome XII
  JOURNAL   Nature 387 (6632 Suppl), 87-90 (1997)
  MEDLINE   97313267
REFERENCE   3  (residues 1 to 397)
  AUTHORS   Saccharomyces Genome Database (yeast-curator@genome.stanford.edu).
  TITLE     Direct Submission
  JOURNAL   Submitted (17-NOV-1999) Department of Genetics, Stanford
            University, Saccharomyces Genome Database, Stanford, CA 94305-5120,
            USA
COMMENT     REFSEQ: This reference sequence was provided by the Saccharomyces
            Genome Database (SGD).
            Method: conceptual translation.
```
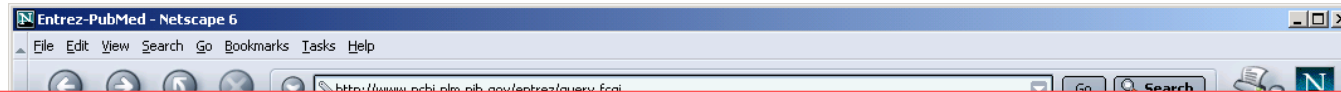
Document: Done (0.721 secs)

NCBI



: J Biol Chem 2000 Jun 23;275(25):19352-60                    Related Articles,  NEW  Links

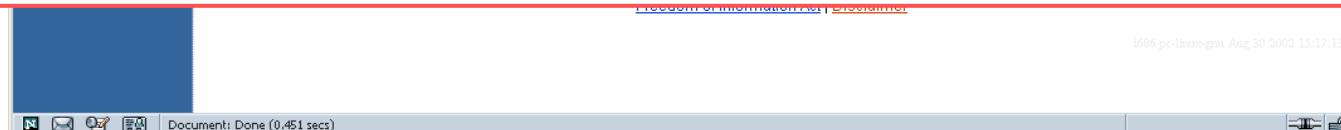FREE full text article at
www.jbc.org

## Psr1p/Psr2p, two plasma membrane phosphatases with an essential DXDX(T/V) motif required for sodium stress response in yeast.

Siniossoglou S, Hurt EC, Pelham HR.

Medical Research Council Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, United Kingdom.

Regulation of intracellular ion concentration is an essential function of all cells. In this study, we report the identification of two previously uncharacterized genes, PSR1 and PSR2, that perform an essential function under conditions of sodium ion stress in the yeast Saccharomyces cerevisiae. Psr1p and Psr2p are highly homologous and were identified through their homology with the endoplasmic reticulum membrane protein Nem1p. Localization and biochemical fractionation studies show that Psr1p is associated with the plasma membrane via a short amino-terminal sequence also present in Psr2p. Growth of the psr1psr2 mutant is severely inhibited under conditions of sodium but not potassium ion or sorbitol stress. This growth defect is due to the inability of the psr1psr2 mutant to properly induce transcription of ENA1/PMR2, the major sodium extrusion pump of yeast cells. We provide genetic evidence that this regulation is independent of the phosphatase calcineurin, previously implicated in the sodium stress response in yeast. We show that Psr1p contains a DXDX(T/V) phosphatase motif essential for its function in vivo and that a Psr1p-PtA fusion purified from yeast extracts exhibits phosphatase activity. Based on these data, we suggest that Psr1p/Psr2p, members of an emerging class of eukaryotic phosphatases, are novel regulators of salt stress response in yeast.

PMID: 10777497 [PubMed - indexed for MEDLINE]

# New Entrez Databases

## 6 new databases in the last year

1. Books: online books
2. GEO: high-throughput gene expression and microarray datasets
3. 3D Domains: structural protein domains from Entrez Structure
4. UniSTS: markers and mapping data
5. CDD: conserved protein domains
6. SNP: single nucleotide polymorphisms

## 5 new databases on the way

1. UniGene: clusters of sequence similar transcripts
2. Gene: a derivation of LocusLink and Genomes
3. SKY/CGH: spectral karyotyping/comparative genomic hybridization
4. Site Search: search the NCBI web and ftp sites
5. Gensat: *in situ* gene expression in the nervous system of the mouse

# Entrez Gensat

# Current Query Scheme

Database selection

Query

Records

links

# Global Query Scheme

NCBI

# Entrez Global Query

# National Center for Biotechnology Information
National Library of Medicine    National Institutes of Health

NCBI

| PubMed | Entrez | BLAST | OMIM | Books | TaxBrowser | Structure |

Search [Nucleotide ▼] for [          ] [Go]

▶ What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease.

Try these: Map

### The new buzz on UniGene!

The fly and mosquito join the UniGene sequences collection, an experimental system for automatically partitioning GenBank sequences into a non-redundant set of gene-oriented clusters. New additions serve to further enhance UniGene as a public resource for gene discovery. More...

▶ NCBI in the News

NCBI now offers quick links to online resources through LinkOut. This new feature of the Entrez database system "expands the biological relevance of NCBI's molecular

## Hot Spots

▶ Cancer genome anatomy project

▶ Clusters of orthologous groups

▶ Coffee Break

▶ Electronic

▶ Human map view

▶ Human/mouse homology maps

▶ LocusLink

▶ Malaria genetics & genomics

▶ Mouse genome resources

▶ ORF finder

▶ Reference sequence project

# www.ncbi.nlm.nih.gov